Language
Technologies
Institute

Carnegie
Mellon
University

# Multimodal Machine Learning
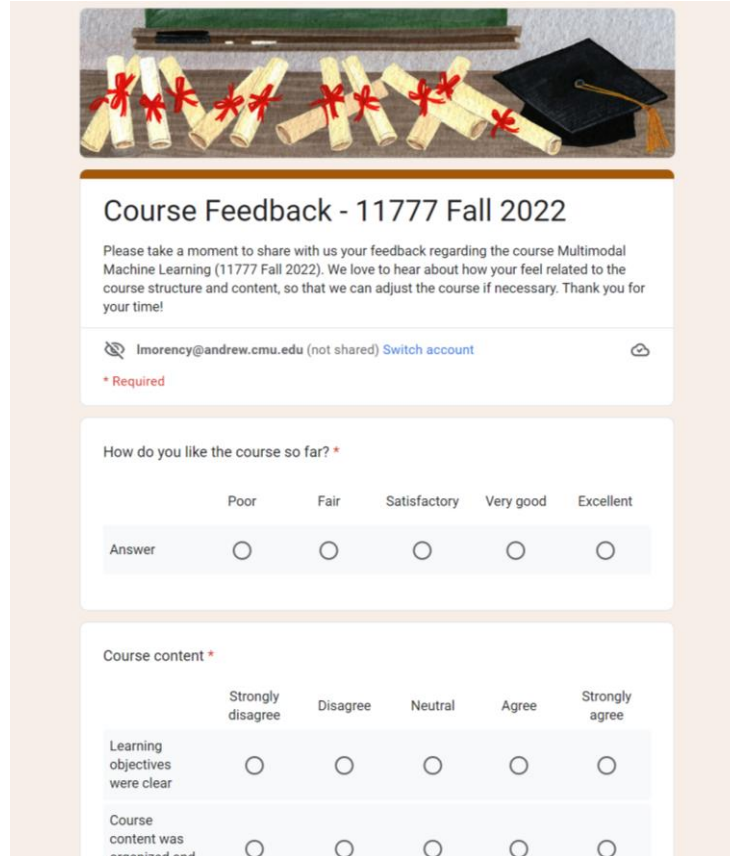
## Lecture 5.2: Aligned Representations

Louis-Philippe Morency

*\* Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yonatan Bisk.*

# Administrative Stuff

# Share Your Thoughts!

https://forms.gle/8vmWa7PxBfkGC2i69



**Deadline**
Please submit your feedback about this course before this Wednesday 10/5

Optional,
but greatly appreciated! ☺

Anonymous, by default.
- You can optionally share your email address if you want us to follow-up with you directly.

# Second Project Assignment (Due Monday 10/10)

Main goals:

1. Help clarify and expand your research ideas
   - Build qualitative intuitions by directly studying the original data
   - Perform analyses on your dataset, relevant to your research ideas
2. Understand the structure in your data and modalities
   - Perform analyses and visualizations to understand each modality
   - Study representations from CNNs, word2vec, BERT, …

Two types of analyses:
- Idea-oriented analyses
- Modality-oriented analyses

# Second Project Assignment (Due Monday 10/10)

Examples of **idea-oriented** analyses:

- What external knowledge is needed when performing the task?
- How often multimodal information is needed? How is it integrated?
- What biases may be present in the data? Which modalities?

Examples of **modality-oriented** analyses:

- What are the different verbs used in the VQA questions?
- What objects do not get detected? Are they important?
- Visualize face embeddings with respect of emotion labels

Language Technologies Institute

Carnegie Mellon University

# Second Project Assignment (Due Monday 10/10)

Idea-oriented analyses:

- **Human simulations:** Instead of a computer, try to do the same task as a human. Gather notes on how you perform the task.

- **Data analysis:** study the multimodal data (e.g., using statistical methods) to clarify your hypotheses related to your research ideas

Modality oriented analyses:

- **Language modality:** explore the language structure in your dataset. You can compare word-level and sentence-level embeddings.

- **Visual modality:** study visual representations for your dataset. You visualize how your visual features successfully model your labels.

# Second Project Assignment (Due Monday 10/10)

Number of analyses:

- Teams of 3 or 4 students: 2 analyses (4 pages)
- Teams of 5 or 6 students: 3 analyses (6 pages)


➢ You can mix and match between idea-oriented and modality-oriented

➢ Be sure to talk with your TA about formalizing your analysis plan

➢ Each analysis need a separate discussion section


Detailed instructions on Piazza (Resources section)

**Language Technologies Institute**

# Multimodal Machine Learning

## Lecture 5.2: Aligned Representations

**Louis-Philippe Morency**

# Objectives of today's class

- Contextualized sequence representations
- Transformer networks
    - Self-attention
    - Multi-head attention
    - Position embeddings
    - Sequence-to-sequence modeling
- Multimodal contextualized embeddings
- Language pre-training
    - BERT pre-training and fine-tuning

# Contextualized Sequence Representations

# Sequence Encoding - Contextualization

**Option 1: Bi-directional LSTM:**
(e.g., ELMO)



How to encode this sequence while modeling the interaction between elements (e.g., words)?

But harder to parallelize…

# Sequence Encoding - Contextualization

**Option 2: Convolutions**



Can be parallelized!

But modeling long-range dependencies
require multiple layers

And convolutional kernels are static

# Sequence Encoding - Contextualization

**Option 3: Self-attention**

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

Contextualized Sequence Encoding

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

I　do　not　like　it

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

**Self-Attention**

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

I　do　not　like　it

Can be parallelized!

Long-range dependencies

Dynamic attention weights

# Self-Attention

Carnegie Mellon University

# Self-Attention

# Self-Attention

# Transformer Self-Attention

# Transformer Self-Attention

# Transformer Self-Attention

# Transformer Self-Attention

# Transformer Self-Attention

What if we want to attend simultaneously to multiple subspaces of $x$?



| $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ |

**Transformer's Self-Attention Layer**

$W_q$  $W_k$  $W_v$

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

I    do    not    like    it

# Transformer Multi-Head Self-Attention

# Transformer Multi-Head Self-Attention

# Transformer Multi-Head Self-Attention



$h_1$    $h_2$    $h_3$    $h_4$    $h_5$

## Transformer's Multi-Head Self-Attention Layer

$W_q^3$    $W_k^3$    $W_v^3$

$W_q^2$    $W_k^2$    $W_v^2$

$W_q^1$    $W_k^1$    $W_v^1$

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$

not     like     I     it     do

What happens if the words are shuffled?

# Position embeddings

❑ Position information is not encoded in a self-attention module

How can we encode position information?

**Simple approach:** one-hot encoding

$$
\begin{array}{c}
0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0
\end{array}
\qquad
\begin{array}{c}
0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0
\end{array}
\qquad
\begin{array}{c}
1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0
\end{array}
\qquad
\begin{array}{c}
0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1
\end{array}
\qquad
\begin{array}{c}
0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0
\end{array}
$$

| $x_1$ | $p_1$ | | $x_2$ | $p_2$ | | $x_3$ | $p_3$ | | $x_4$ | $p_4$ | | $x_5$ | $p_5$ |

| not | like | I | it | do |

# Position embeddings

❑ Position information is not encoded in a self-attention module

How can we encode position information?

**Simple approach:** one-hot encoding + linear embeddings + $\begin{cases} \text{Sum} \\ \text{- or -} \\ \text{concat} \end{cases}$



| $x_1$ $p_1$ | $x_2$ $p_2$ | $x_3$ $p_3$ | $x_4$ $p_4$ | $x_5$ $p_5$ |
| not | like | I | it | do |

# Transformer Multi-Head Self-Attention

# Transformer Multi-Head Self-Attention

In vector format…

$h$

Transformer's Multi-Head
Self-Attention Layer

$W_q^3$ $W_k^3$ $W_v^3$

$W_q^2$ $W_k^2$ $W_v^2$

$W_q^1$ $W_k^1$ $W_v^1$

$p$

$x$

# Transformer Multi-Head Attention

# Transformer – Residual Connection

Carnegie Mellon University

# Sequence-to-Sequence Using Transformer

# Sequence-to-Sequence Modeling

Je    n'    aime    pas    cela

$\hat{y}_1$  $\hat{y}_2$  $\hat{y}_3$  $\hat{y}_4$  $\hat{y}_5$

How can we perform seq2seq
translation with transformer attention?

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

I     do    not    like    it

# Seq2Seq with Transformer Attentions

Je        n'        aime     pas      cela

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$   $\hat{y}_4$   $\hat{y}_5$

$h_1$   $h_2$   $h_3$   $h_4$   $h_5$

self-attention

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$

I        do        not       like      it

# Seq2Seq with Transformer Attentions

Je     n'     aime     pas     cela

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$   $\hat{y}_4$   $\hat{y}_5$

| $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ |

**self-attention**

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

I     do     not     like     it

| $g_1$ | $g_2$ | $g_3$ | $g_4$ | $g_5$ |

**"masked" self-attention**

| $y_0$ | $y_1$ | $y_2$ | $y_3$ | $y_4$ |

START     Je     n'     aime     pas

# Seq2Seq with Transformer Attentions

How should we connect the encoder and decoder self-attention to the transformer attention?

Je    n'    aime    pas    cela

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$   $\hat{y}_4$   $\hat{y}_5$

**Transformer attention**

$W_q$    $W_k$    $W_v$

**Query**    **Key**    **Value**

Vector format → $h$

$g$

self-attention

"masked" self-attention

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$

I    do    not    like    it

$y_0$   $y_1$   $y_2$   $y_3$   $y_4$

START    Je    n'    aime    pas

# Seq2Seq with Transformer Attentions

# Language Pre-training

# Token-level and Sentence-level Embeddings

Token-level embeddings

Sentence-level embedding

# Pre-Training and Fine-Tuning



Pre-training
(e.g., language model)

Fine-Tuning

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

① Jointly learn representation for token-level and sentence level

② Same network architecture for pre-training and fine-tuning

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

① Jointly learn representation for token-level and sentence level

② Same network architecture for pre-training and fine-tuning

③ Can be used learn relationship between sentences

④ Models bidirectional and long-range interactions between tokens

How can we do all this?

$h_s$ $h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_{sep}$ $h'_1$ $h'_2$ $h'_3$ $h'_4$ $h'_5$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ sep $x'_1$ $x'_2$ $x'_3$ $x'_4$ $x'_5$

I    do   not   like   it          I   enjoy   my   time   here

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

①     Jointly learn representation for token-level and sentence level

②     Same network architecture for pre-training and fine-tuning

③     Can be used learn relationship between sentences

④     Models bidirectional interactions between tokens

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

**Transformer Self-Attention**

Special sentence-level token

But how to train self-supervised?

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I     do     not     like     it        I     enjoy     my     time     here

# Pre-training BERT Model

**1** **Masked Language Model**

Randomly mask input tokens and then try to predict them

What is the loss function?

# Pre-training BERT Model

② **Next Sentence Prediction**

Given two sentences, predict if this is the next one or not

What is the loss function?

IsNext
**- or -**
NotNext

Where can we find training data?

How can BERT know the difference between both sentences?

$h_s$ $h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_{sep}$ $h'_1$ $h'_2$ $h'_3$ $h'_4$ $h'_5$

## Transformer Self-Attention

cls $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ sep $x'_1$ $x'_2$ $x'_3$ $x'_4$ $x'_5$

I do not like it I enjoy my time here

# Three Embeddings: Token + Position + Sentence



| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

# Fine-Tuning BERT

① Sentence-level classification for only one sentence

Examples: sentiment analysis, document classification

**How?**

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |
|---|---|---|---|---|---|---|---|---|---|---|---|

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  | I | do | not | like | it |  | I | enjoy | my | time | here |

# Fine-Tuning BERT

① Sentence-level classification for only one sentence

Examples: sentiment analysis, document classification

**And if we have a label for each token?**

$\hat{y}_s$

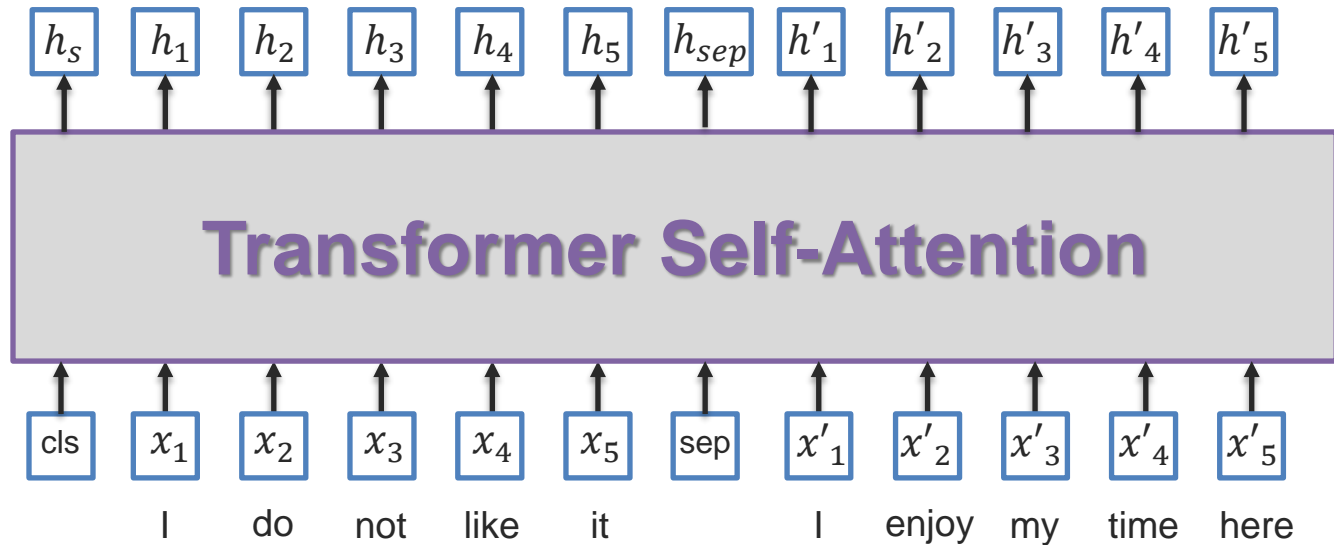| softmax |

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ |

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

I    do    not    like    it

# Fine-Tuning BERT

② Token-level classification for only one sentence

Examples: part-of-speech tagging, slot filling



How to compare two sentences?

# Fine-Tuning BERT

③ Sentence-level classification for two sentences

Examples: natural language inference



$\hat{y}_s$

softmax

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I  do  not  like  it  I  enjoy  my  time  here

# Fine-Tuning BERT

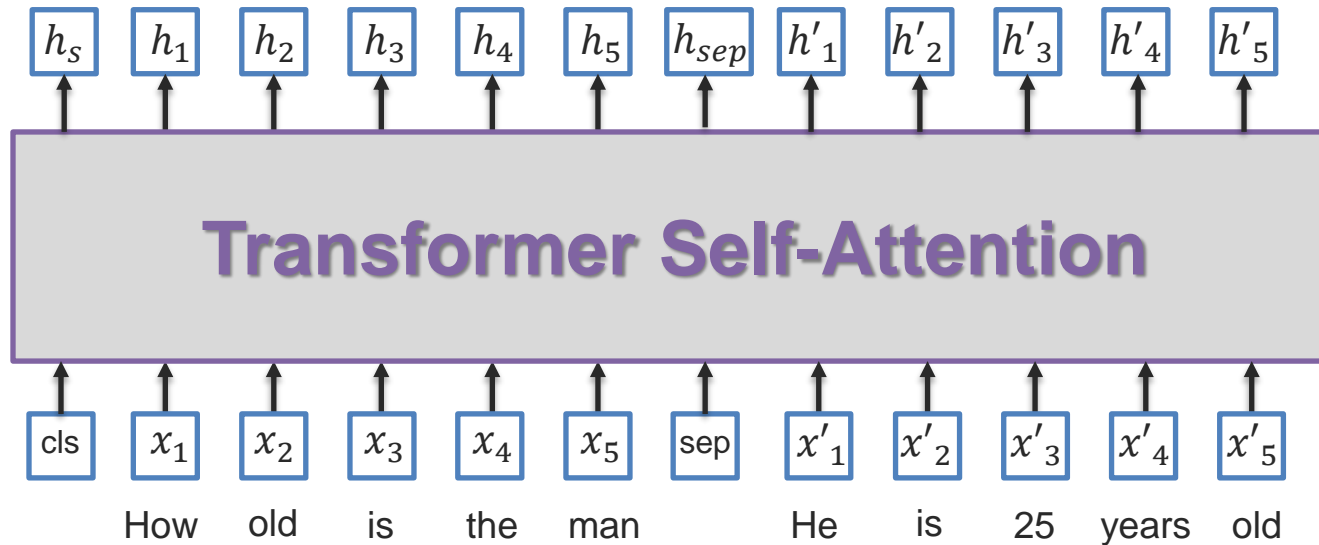④   Question-answering: find start/end of the answer in the document

**Paragraph:** " ... *Other legislation followed, including the Migratory Bird Conservation Act of 1929, a 1937 treaty prohibiting the hunting of right and gray whales, and the Bald Eagle Protection Act of 1940. These later laws had a low cost to society—the species were relatively rare—and little opposition was raised.*"

**Question 1:** "*Which laws faced significant opposition?*"
**Plausible Answer:** *later laws*

**Question 2:** "*What was the name of the 1937 treaty?*"
**Plausible Answer:** *Bald Eagle Protection Act*

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

**Transformer Self-Attention**

How?

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

How   old   is   the   man    He   is   25   years   old

# Fine-Tuning BERT

④ Question-answering: find start/end of the answer in the document



Same architecture for the end time

Maximum value gives start time

softmax

Learned during fine-tuning

$h_s$ $h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_{sep}$ $h'_1$ $h'_2$ $h'_3$ $h'_4$ $h'_5$

**Transformer Self-Attention**

cls $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ sep $x'_1$ $x'_2$ $x'_3$ $x'_4$ $x'_5$

How old is the man    He is 25 years old

# And Many More…     Next week!